

PRESS statistic = sum of squares of
PRESS residuals

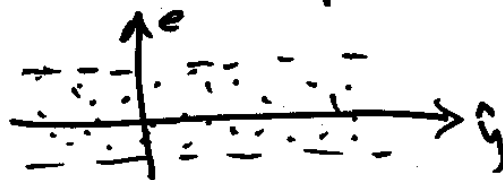
Stat 524
10-24-17
①

$$\text{and } R^2_{\text{PRESS}} = 1 - \frac{\text{PRESS stat.}}{SST}$$

(compare to $R^2 = 1 - \frac{SSE}{SST}$)

Residual Diagnostics

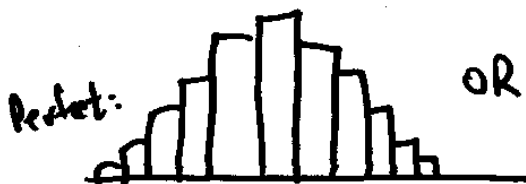
① Plot the residuals vs. the predicted values



Perfect.
- no change in
dispersion
- no pattern

② Create a histogram or a box plot
of the residuals

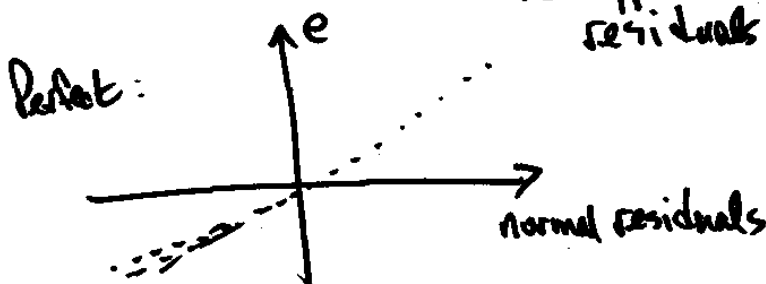
②



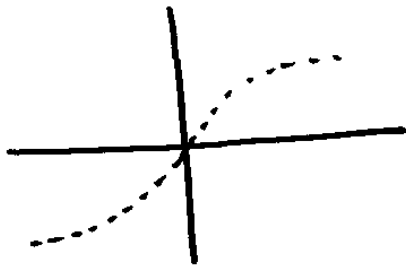
OR



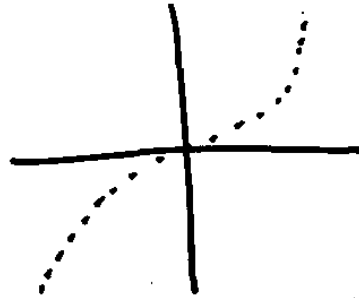
③ Q-Q plot: Plot the residuals vs.
the hypothesized normally distributed
residuals



③

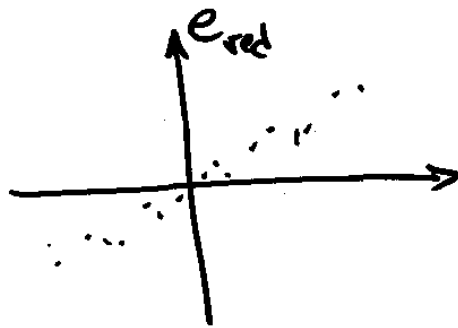


Your residuals have
shorter tails
than normal



Your residuals have
longer tails than
normal

- ④ Remove 1 predictor from the model.
Run the reduced model + obtain the residuals.
Plot the residuals vs. the removed predictor.

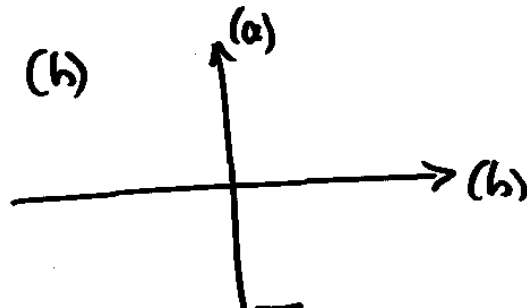


④ If you can see a
pattern, then X_i
should not have
been removed from
the model

If there is no discernible pattern, then
 X_i is a candidate for removal from the model.

- ⑤ Remove 1 predictor from the model.
Run the reduced model & obtain the residuals (a)
Regress the removed predictor on the remaining predictors
+ obtain those residuals (b)

Plot (a) vs (b)



(5)

Same interpretation as in 14

16 Keep track of order of data collection

$$\text{let } r_k = \frac{\sum_{t=k+1}^n e_t e_{t-k}}{\sum_{t=1}^n e_t^2} \quad \text{for } k=1, 2, \dots, n-1$$

These are the serial autocorrelations

(6)

Ideally, these are all 0.

Fact: Under the assumption of independent errors,

$$E[r_k] = 0 \quad \text{and} \quad V[r_k] \approx \frac{1}{n}$$

So compute $\frac{r_k}{\sqrt{\frac{1}{n}}}$ + find any that are outside of ± 2

[7] Durbin-Watson test

(7)

$$D = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}$$

(related to the
k=1 case
of [6])

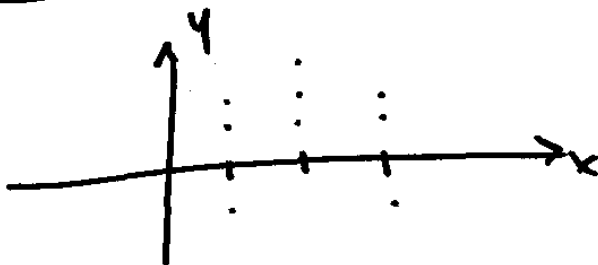
Under the assumption of independent errors,

$$E[D] = 2$$

If the DW test is significant,
you have found a violation of your assumption.

"Lack of fit" Test

(8)



Must have multiple observations of y
at at least 1 x-value

For simple linear regression, assume we have n
observations, but m distinct x-values

π_1 : observe y_{11}, \dots, y_{1n_1}

π_2 : y_{21}, \dots, y_{2n_2}

\vdots

π_m : y_{m1}, \dots, y_{m,n_m}

(9)

$$n_1 + n_2 + \dots + n_m = n$$

$$SSE = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \hat{y}_i)^2$$

$$= \sum_{i=1}^m \sum_{j=1}^{n_i} (\underbrace{y_{ij} - \bar{y}_i}_{(1)} + \underbrace{\bar{y}_i - \hat{y}_i}_{(2)})^2$$

$$= \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2$$

(1) (2)

$$+ 2 \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)(\bar{y}_i - \hat{y}_i)$$

(3)

$$\textcircled{3}: 2 \sum_{i=1}^m \left[(\bar{y}_i - \hat{y}_i) \underbrace{\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)}_{\sum_{j=1}^{n_i} y_{ij} - n_i \bar{y}_i} \right]$$

$$\begin{aligned}
 \text{So } SSE = & \underbrace{\sum_{i=1}^m \sum_{j=1}^{n_j} (y_{ij} - \bar{y}_i)^2}_{\text{"pure error"} \quad SS_{PE}} + \underbrace{\sum_{i=1}^m n_j (\bar{y}_i - \hat{\bar{y}}_i)^2}_{\text{"lack of fit"} \quad SS_{LF}} \quad (11)
 \end{aligned}$$

The lack of fit F test is

$$F = \frac{SS_{LF}/(m-2)}{SS_{PE}/(n-m)} \sim F_{m-2, n-m}$$

If this is significant, you have a poor fit.

Note about the df:

$$\begin{aligned}
 (m-2) + (n-m) &= n-2, \text{ which was} \\
 &\text{df in simple linear regression}
 \end{aligned}$$

For multiple regression, $df_{LF} = m-p$

$$df_{PE} = n-m$$

$$\text{So } (m-p) + (n-m) = n-p = df_E$$